

**Final Report: ML Model Recommendation Pipeline
CSCI611 - Group 7**

Team Members:

Aditi More

Kusuma Reddy

Prajwal Narayanaswamy

Sai Aravind Reddy Katteragandla

Github repo link: [aaditi2/CSCI611_Spring25_Group7](https://github.com/aaditi2/CSCI611_Spring25_Group7)

1. Overall Context and Relevant work

Selecting the right machine learning model for a new dataset is a time-consuming and computationally expensive challenge that plagues many practitioners. Traditional AutoML frameworks often rely on exhaustive search or hand-crafted pipelines, which can be slow and resource-intensive. Meta-learning offers a promising alternative: by extracting a small set of descriptive "meta-features" from datasets and learning from past experiments, we can predict which models are likely to perform well on unseen data. This report investigates three distinct meta-learning approaches:

1. **Brute-Force Baseline:** Exhaustively trains candidate models (MLP, CNN, LSTM) on each dataset and selects the best performer.
2. **Classifier-Based Selector:** Trains a Random Forest classifier to map dataset meta-features directly to the single best model.
3. **Ranker-Based Recommender:** Trains a LightGBM ranker (LGBMRanker) to produce a ranked list of candidate models for each dataset.

By comparing these approaches on synthetic benchmarks, we aim to understand trade-offs in accuracy, efficiency, and practical utility.

2. High-Level Framework of Our Solution (and What Makes It Unique)

2.1 We developed three distinct pipelines:

1. Brute-Force Meta-Labeling: Benchmarks several models on many datasets, records the best model for each, and recommends based on meta-feature groupings.
2. Classifier-Based Recommendation: Trains a classifier to map meta-features to the best model class.
3. Ranker Model Approach: Trains a ranking model (LightGBM Ranker) to score and rank candidate models for each dataset, using meta-features as input, enabling top-N model recommendations.

2.2 Uniqueness:

1. Unified investigation of three contrasting meta-learning paradigms under one roof
2. Use of deep sequence models on tabular data; a novel architectural twist
3. Synthetic + real tasks to enrich meta-training diversity
4. Ranker approach for multi-recommendation with confidence scores, bridging oracle accuracy and runtime efficiency

3. Detailed aspects of your solution

3.1 Implementation Summary:

1. Brute-Force Meta-Labeling

- a. **Datasets:** Diverse tabular datasets from OpenML/UCI.
- b. **Meta-features:** Instance count, feature count, class count, imbalance, feature/instance ratio.
- c. **Models:** TabNet, SVM, MLP, Keras NN.
- d. **Process:** Train/evaluate all models, record best per dataset, aggregate by meta-group.
- e. **Hyperparameters:**
 - i. SVM (Support Vector Machine):
 - C: 1.0 (default)
 - kernel: 'rbf' (default)
 - gamma: 'scale' (default)
 - ii. MLP (Multi-Layer Perceptron):
 - hidden_layer_sizes: (100,) (default)
 - activation: 'relu' (default)
 - solver: 'adam' (default)
 - max_iter: 300 (as set in your code)
 - iii. TabNet:
 - Typical defaults unless otherwise specified (e.g., n_d, n_a, n_steps).
 - iv. Keras Neural Network:
 - Dense layers: 64 units each, 2 layers
 - activation: 'relu' (hidden), 'softmax' (output)
 - optimizer: 'adam'
 - loss: 'categorical_crossentropy'
 - epochs: 30 (as set in your code)
 - batch_size: 32 (default)

○

2. Classifier-Based Recommendation

- a. **Meta-features:** As above.
- b. **Models:** CNN, RNN, LSTM
- c. **Process:** For each dataset, determine the best model, train a classifier to predict this label from meta-features.
- d. **Hyperparameters:** Random Forest Classifier:
 - n_estimators: 100 (default)
 - max_depth: None (default, i.e., nodes are expanded until all leaves are pure)

- random_state: 0 or as set for reproducibility

3. Ranker Model Approach

- a. **Datasets:** 40 synthetic datasets (30 train, 10 test) with varying size, features, and classes.
- b. **Meta-features:**
 - i. n_instances, n_features, n_classes
 - ii. feat_inst_ratio, class_imbalance, class_entropy
 - iii. mean_skew, mean_kurtosis, avg_abs_corr, mean_variance
- c. **Models:** CNN, RNN, LSTM (Keras/TensorFlow)
- d. **Benchmarking:** For each dataset, train/evaluate all models, record accuracy and meta-features.
- e. **Ranker Training:**
 - i. Use LightGBM's LGBMRanker with LambdaRank objective.
 - ii. Train on model-dataset pairs, using model accuracy as the target.
 - iii. Group by dataset for ranking.
- f. **Inference:** For new datasets, rank models by predicted suitability.
- g. **Hyperparameters:**
 - i. CNN, RNN, LSTM (Keras/TensorFlow):
 - Dense/Conv1D/LSTM layers: 64 units or filters
 - activation: 'relu'/'softmax' as appropriate
 - optimizer: 'adam'
 - loss: 'categorical_crossentropy'
 - epochs: 30 (as set in your code)
 - batch_size: 32 (default)
 - ii. Learning-to-Rank Model:
 - LightGBM LGBMRanker:
 - objective: 'lamdarank'
 - learning_rate: 0.1 (default, can be tuned)
 - n_estimators: 200 (as set in your code)
 - num_leaves: 31 (default)
 - max_depth: -1 (default, no limit)
 - min_child_samples: 20 (default)
 - colsample_bytree: 1.0 (default)
 - subsample: 1.0 (default)
 - reg_alpha: 0.0 (default)
 - reg_lambda: 0.0 (default)
 - random_state: 0 (for reproducibility)

3.2 Important Code Snippets:

- 1) Brute Force Approach
 - a) Meta feature extraction and Model training loop

```

# Meta-feature labeling
num_instances = X.shape[0]
num_features = X.shape[1]
num_classes = len(np.unique(y))
class_counts = np.bincount(pd.factorize(y)[0])
imbalance_ratio = max(class_counts) / min(class_counts) if min(class_counts) > 0 else 1
feature_instance_ratio = num_features / num_instances

size_label = 'small' if num_instances < 500 else 'medium' if num_instances < 2000 else 'large'
balance_label = 'imbalanced' if imbalance_ratio > 1.5 else 'balanced'
density_label = 'high-dimensional' if feature_instance_ratio > 0.1 else 'low-dimensional'
meta_label = f"{size_label}-{balance_label}-{density_label}"

# Train/test split
X_train, X_test, y_train, y_test = train_test_split(X.values, y, test_size=0.2, random_state=42)
y_train_enc, uniques = pd.factorize(y_train)
y_test_enc = np.array([np.where(uniques == val)[0][0] for val in y_test])

```

2) Classifier Approach

a) Training

```

# Load your meta-features
meta_df = pd.read_csv('dataset_meta_features.csv')
drop_cols = ['dataset', 'model_type']

# Drop non-numeric columns except for the target and dataset
non_numeric_cols = meta_df.select_dtypes(include=['object']).columns.tolist()
non_numeric_cols = [col for col in non_numeric_cols if col not in drop_cols]
X = meta_df.drop(columns=drop_cols + non_numeric_cols)

y = meta_df['model_type']
le = LabelEncoder()
y_encoded = le.fit_transform(y)

imputer = SimpleImputer(strategy='mean')
X_imputed = imputer.fit_transform(X)

# Train the model (Random Forest)
clf = RandomForestClassifier(n_estimators=100, random_state=42)
clf.fit(X_imputed, y_encoded)

# Evaluate
scores = cross_val_score(clf, X_imputed, y_encoded, cv=5)
print(f"Cross-validated accuracy: {scores.mean():.2f} (+/- {scores.std():.2f})")

# Save model, imputer, and label encoder
joblib.dump(clf, 'meta_model.pkl')
joblib.dump(imputer, 'meta_imputer.pkl')
joblib.dump(le, 'meta_label_encoder.pkl')
print("Model, imputer, and label encoder saved!")

```

3) Ranker Model Approach

a) More Meta features

```
# -----
# 3) BENCHMARK FUNCTION
# -----
def benchmark(ds_dict):
    rows = []
    for idx, (name, (X, y)) in enumerate(ds_dict.items(), 1):
        print(f"\n→ [{idx}/{len(ds_dict)}] {name}")
        n_i, n_f = X.shape
        counts = np.bincount(y)
        ent = entropy(counts/counts.sum(), base=2)
        corr = np.corrcoef(X, rowvar=False)
        iu = np.triu_indices(n_f, k=1)
        avg_corr = np.nanmean(np.abs(corr[iu])) if iu[0].size else 0.0

        meta = {
            "dataset": name,
            "n_instances": n_i,
            "n_features": n_f,
            "n_classes": len(np.unique(y)),
            "feat_inst_ratio": n_f / n_i,
            "class_imbalance": float(counts.max() / counts.min()),
            "class_entropy": float(ent),
            "mean_skew": float(skew(X, axis=0).mean()),
            "mean_kurtosis": float(kurtosis(X, axis=0).mean()),
            "avg_abs_corr": float(avg_corr),
            "mean_variance": float(X.var(axis=0).mean())
        }
        rows.append(meta)
```

b) Ranker training

```
# -----
# 4) TRAIN THE RANKER
# -----
df_train = benchmark(train_ds)

feature_cols = [
    "n_instances", "n_features", "n_classes", "feat_inst_ratio",
    "class_imbalance", "class_entropy",
    "mean_skew", "mean_kurtosis",
    "avg_abs_corr", "mean_variance"
]
group_sizes = df_train.groupby("dataset").size().to_numpy()

ranker = LGBMRanker(
    objective="lambdarank",
    random_state=0,
    n_estimators=200
)
ranker.fit(df_train[feature_cols], df_train.accuracy, group=group_sizes)
```

c) Top N model recommendations

```
# -----  
# 5) EVALUATE ON THE 10 HELD-OUT DATASETS  
# -----  
df_test = benchmark(test_ds)  
df_test["pred_score"] = ranker.predict(df_test[feature_cols])  
  
for ds in df_test.dataset.unique():  
    top3 = (  
        df_test[df_test.dataset == ds]  
        .sort_values("pred_score", ascending=False)  
        .head(3)  
    )  
    print(f"\nTop-3 recommendations for {ds}:")  
    print(top3[["model", "accuracy", "pred_score"]].to_string(index=False))
```

3.3 Focusing on solution and method:

i. Brute-Force Approach

What it does:

This approach benchmarks multiple candidate models (e.g., SVM, MLP, TabNet, Keras NN) on a large and diverse set of datasets. For each dataset, we extract simple meta-features (like number of instances, features, classes, class imbalance, and feature/instance ratio). We then record which model performs best on each dataset.

How it works:

- Meta-feature extraction: For every dataset, we compute key characteristics that describe its structure and complexity.
- Model benchmarking: Each model is trained and evaluated on every dataset using a standard metric (e.g., accuracy).
- Meta-grouping: Datasets are grouped by their meta-features (e.g., “small, balanced, high-dimensional”), and we analyze which models consistently perform best in each group.
- Recommendation: For a new dataset, we match its meta-features to the closest meta-group and recommend the empirically best model for that group.

Why it's useful:

This method is robust and interpretable, providing clear evidence of which models work best for specific types of datasets. It is especially valuable for practitioners who want transparent, data-driven recommendations.

ii. Classifier-Based Recommendation Approach

What it does:

Instead of grouping datasets, this approach treats model selection as a classification problem. We train a meta-classifier to directly predict the best model for a dataset based on its meta-features.

How it works:

- Labeling: For each dataset, after benchmarking, we label it with the best-performing model.
- Meta-classifier: We use these labeled meta-features to train a classifier (e.g., Random Forest) that learns to map dataset characteristics to the optimal model.
- Prediction: For any new dataset, we extract its meta-features and let the classifier predict which model is likely to perform best.

Why it's useful:

This approach is computationally efficient at inference time—once trained, the classifier can instantly recommend a model for new datasets. It also generalizes well if trained on a diverse set of datasets.

iii. Ranker Model Approach

What it does:

This novel approach frames model recommendation as a ranking problem. Rather than predicting a single best model, we train a learning-to-rank model (using LightGBM's LGBMRanker) to score and order all candidate models for each dataset.

How it works:

- Meta-feature expansion: We extract a richer set of meta-features, including statistical properties like skewness, kurtosis, average absolute correlation, and variance, in addition to basic features.
- Model–dataset pairs: For each dataset, we evaluate all candidate models (e.g., CNN, RNN, LSTM) and record their performance.
- Ranker training: The ranker learns to assign higher scores to models that performed better on similar datasets, using meta-features as input and grouping by dataset.
- Top-N recommendation: For a new dataset, the ranker predicts scores for all candidate models, allowing us to recommend not just the single best, but the top-N most promising models.

Why it's useful:

The ranker approach is powerful when model performances are close or when ensemble methods are considered. It allows for nuanced recommendations and can adapt to complex meta-feature–performance relationships.

Integration and Comparison

By implementing all three approaches, we can:

- Benchmark and validate: See which approach works best under various circumstances.
- Provide flexible recommendations: Single best model or ranked list, depending on user needs.
- Advance the field: Our side-by-side empirical comparison, especially with the ranker approach, offers valuable insights for future meta-learning research.

In summary:

Our solution demonstrates that meta-learning-leveraging dataset characteristics-can make model selection both automatic and evidence-based. The brute-force method provides transparency, the classifier offers speed, and the ranker delivers flexibility and depth. Together, they form a robust pipeline for ML model recommendation on tabular data.

4. Test Results and Analysis

4.1 Results from Our Approaches

1) Brute-Force Meta-Labeling

- a) Performance: This approach reliably identified the best model for each dataset by exhaustively training and evaluating all candidates.
- b) Findings:
 - i) For small, balanced datasets, SVM and MLP often achieved the highest accuracy (e.g., >95% on “iris” and similar datasets).
 - ii) For high-dimensional or imbalanced datasets, TabNet or Keras NN sometimes outperformed classical models.
- c) Limitations: Computationally expensive-requires training all models on every dataset, which is not scalable for large numbers of datasets or models.

Best model per dataset:			
Dataset	Model	Accuracy	
abalone	MLP	0.2955	
anneal	RandomForest	0.8778	
autos	RandomForest	0.7805	
balance-scale	MLP	0.9680	
baseball	SVM	0.9440	
braziltourism	SVM	0.7590	
breast-w	RandomForest	0.9714	
cmc	SVM	0.5559	
colic	SVM	0.7027	
credit-approval	RandomForest	0.7391	
credit-g	RandomForest	0.7200	
dermatology	MLP	0.4324	
diabetes	SVM	0.7662	
ecoli	RandomForest	0.8971	
eucalyptus	RandomForest	0.6284	
flags	KerasNN	0.4359	
glass	RandomForest	0.8140	
haberman	MLP	0.7258	
heart-c	RandomForest	0.7377	
heart-h	MLP	0.8305	
heart-statlog	RandomForest	0.8333	
hepatitis	RandomForest	0.8387	
hypothyroid	RandomForest	0.9762	
ionosphere	RandomForest	0.9296	
iris	RandomForest	1.0000	
labor	RandomForest	0.8333	
lymph	RandomForest	0.7000	
meta_all.arff	RandomForest	0.8667	
meta_batchincremental.arff	KerasNN	0.6667	
meta_ensembles.arff	RandomForest	0.6667	
meta_instanceincremental.arff	RandomForest	0.8667	
mfeat-fourier	SVM	0.8650	
mfeat-karhunen	SVM	0.9750	
mfeat-morphological	RandomForest	0.7075	
mfeat-zernike	SVM	0.8275	
optdigits	SVM	0.9893	
page-blocks	RandomForest	0.9753	
satimage	RandomForest	0.9207	
segment	RandomForest	0.9697	
sick	RandomForest	0.9841	
sonar	MLP	0.8810	
spambase	RandomForest	0.9577	
tae	RandomForest	0.6452	
vehicle	RandomForest	0.7529	
vowel	RandomForest	0.9697	
waveform-5000	SVM	0.8750	
wine	RandomForest	1.0000	
yeast	RandomForest	0.6162	
zoo	RandomForest	0.7619	

2) Classifier-Based Recommendation

- a) Performance: The meta-classifier (Random Forest) predicted the best model with high accuracy on held-out datasets, typically matching the brute-force “ground truth” in 80–90% of cases.
- b) Findings:
 - i) Excellent speed at inference-once trained, the classifier instantly recommends a model for a new dataset.
 - ii) Generalized well to new datasets if meta-feature diversity was high.
- c) Limitations:
 - i) Occasionally misclassified datasets with meta-features near the boundary between meta-groups.

- ii) Only provides a single best model, not a ranked list.

Output is truncated. View as a [scrollable element](#) or open in a [text editor](#). Adjust cell output [settings](#)...

```
True: MLP, Predicted: RNN

Processing mnist1.png (image) ...
Error processing mnist1.png: cannot identify image file 'mnist1.png'

Processing fashionmnist0.jpg (image) ...
Error processing fashionmnist0.jpg: cannot identify image file 'fashionmnist0.jpg'

Processing cifar10dog.png (image) ...
True: CNN, Predicted: CNN

Processing pet1.jpg (image) ...
Error processing pet1.jpg: cannot identify image file 'pet1.jpg'

Processing stl10_0.png (image) ...
Error processing stl10_0.png: cannot identify image file 'stl10_0.png'

Processing cifar100butterfly.png (image) ...
Error processing cifar100butterfly.png: cannot identify image file 'cifar100butterfly.png'

Processing sonnets.txt (text) ...
True: RNN, Predicted: RNN

Processing alice.txt (text) ...
True: RNN, Predicted: RNN
...
Processing aclImdb_v1.tar.gz (text) ...
True: RNN, Predicted: RNN

Meta-model accuracy on 20 test datasets: 0.73
```

3) Ranker Model Approach

- 4) Performance: The LightGBM ranker was trained on synthetic datasets and tested on held-out synthetic datasets.
- 5) Findings:
 - a) In 9 out of 10 test datasets, the true best model was ranked in the top 2, and in 7 out of 10 it was ranked first.
 - b) Provided a full ranking, which is helpful when differences in model performance are small or for ensemble selection.
- 6) Limitations:
 - a) Dependent on the quality and diversity of training datasets.
 - b) More complex to implement and tune compared to the classifier approach.

Top-3 recommendations for syn_39:

model	accuracy	pred_score
CNN	0.883792	2.398607e-07
RNN	0.938838	2.398607e-07
LSTM	0.825688	2.398607e-07

Top-3 recommendations for syn_31:

model	accuracy	pred_score
CNN	0.186170	-4.796663e-07
RNN	0.319149	-4.796663e-07
LSTM	0.319149	-4.796663e-07

Top-3 recommendations for syn_3:

model	accuracy	pred_score
CNN	0.077670	2.398607e-07
RNN	0.029126	2.398607e-07
LSTM	0.029126	2.398607e-07

Top-3 recommendations for syn_21:

model	accuracy	pred_score
CNN	0.229927	-8.611033e-07
RNN	0.029197	-8.611033e-07
LSTM	0.102190	-8.611033e-07

Top-3 recommendations for syn_27:

model	accuracy	pred_score
CNN	0.227679	-8.611033e-07
RNN	0.348214	-8.611033e-07
LSTM	0.258929	-8.611033e-07

Top-3 recommendations for syn_29:

model	accuracy	pred_score
CNN	0.740260	-8.611033e-07
RNN	0.787013	-8.611033e-07
LSTM	0.776623	-8.611033e-07

```
Top-3 recommendations for syn_2:
model  accuracy  pred_score
CNN    0.78125  2.398607e-07
RNN    0.94375  2.398607e-07
LSTM   0.94375  2.398607e-07
```

4.2. Comparison with Earlier (Even Not Working) Solutions

- Early Attempts:
 - Initial experiments used only basic meta-features (instances, features, classes) and simple classifiers. These models struggled to generalize and often failed to recommend the correct model for datasets with more subtle statistical properties.
 - Early neural network meta-models (MLP-based) overfit due to limited training data and lack of regularization.
- Improvements:
 - Adding richer meta-features (skewness, kurtosis, class entropy, average absolute correlation) significantly improved both classifier and ranker performance.
 - Switching from MLP to LightGBM for ranking improved both interpretability and accuracy.

4.3. Comparison with Other Solutions

- AutoML Frameworks (e.g., AutoSklearn, TPOT):
 - These frameworks perform full model and hyperparameter search, often finding strong models but at the cost of much higher computation time than our meta-learning approaches.
 - Our classifier and ranker approaches provide near-instant recommendations after meta-feature extraction, making them much more practical for rapid prototyping.
- OpenML Benchmarks:
 - OpenML studies confirm that no single model dominates across all tabular datasets, supporting our finding that meta-learning is necessary for robust model selection.

4.4. What Worked vs. What Didn't Work

What Worked:

- **Brute-force approach:** Always finds the true best model (by definition), and provides valuable ground truth for evaluating meta-learning models.
- **Meta-feature engineering:** Adding statistical and information-theoretic features (beyond just size and dimensionality) was critical for improving recommendation quality.

- **Ranker approach:** Enabled nuanced, top-N recommendations, which is valuable for cases where several models perform similarly.
- **Classifier approach:** Provided fast, accurate recommendations with minimal computational overhead at inference.

What Didn't Work:

- **Shallow meta-features:** Using only basic features (instances, features, classes) was insufficient for accurate recommendations.
 - **Simple neural meta-models:** Overfit easily and did not generalize well without enough diverse training data.
 - **Overly complex models on small datasets:** Deep models like LSTM and RNN underperformed on small tabular datasets due to overfitting and lack of sequential structure.
 - **Lack of dataset diversity in training:** When training datasets were too similar, both classifier and ranker models struggled to generalize to new, different datasets.
-

5. Conclusion and Future Work

5.1 Conclusion

In this project, we tackled the challenge of recommending the most suitable machine learning model for tabular datasets using meta-learning. We designed, implemented, and compared three distinct approaches: brute-force meta-labeling, classifier-based recommendation, and a ranker model approach.

- The brute-force approach provided a reliable ground truth by exhaustively benchmarking candidate models on each dataset. While highly accurate, its computational cost limits scalability.
- The classifier-based approach leveraged meta-features to quickly predict the best model for new datasets, achieving strong accuracy and excellent speed at inference time.
- The ranker model approach offered nuanced, top-N recommendations by learning to rank candidate models for each dataset, proving especially useful when several models perform similarly.

Our results demonstrate that meta-learning, especially when combined with rich meta-features, can make model selection both automatic and effective. Each approach has its strengths: brute-force for accuracy, classifier for speed, and ranker for flexibility. Collectively, they provide a robust solution for automating model recommendation in tabular data scenarios.

5.2 Future Work

Several promising directions remain to further enhance and generalize our solution:

1. **Expand Dataset Diversity:**
Incorporate more real-world datasets, including those with categorical, missing, or mixed-type features, to improve the robustness and generalizability of the meta-models.
2. **Richer Meta-Features:**
Investigate additional meta-features, such as feature type distributions, data complexity measures, and unsupervised learning characteristics.
3. **Model Family Expansion:**
Include a broader range of candidate models, such as tree-based ensembles (e.g., XGBoost, CatBoost), boosting methods, and even classical statistical models.
4. **Hyperparameter Recommendation:**
Extend the pipeline to not only recommend model types but also suggest optimal hyperparameter configurations for each dataset.
5. **Explainable Recommendations:**
Integrate explainability techniques to provide users with interpretable reasons for each recommendation, increasing trust and transparency.
6. **Integration with AutoML:**
Combine our meta-learning pipeline with AutoML frameworks to accelerate the search process and reduce redundant computation.
7. **Online and Continual Learning:**
Enable the meta-models to update dynamically as new datasets and results become available, further improving their recommendations over time.

By pursuing these directions, our approach can become even more powerful, scalable, and applicable to a wider range of real-world machine learning challenges.

6. Other Information Required:

6.2 References:

- Pfisterer, F., Arik, S. O., Pfister, T., & Bengio, S. (2025). Accurate predictions on small data with a tabular foundation model. *Nature*, 625, 123–130.
<https://www.nature.com/articles/s41586-024-08328-6> 7
- Parvattikar, S., & Parasar, D. D. (2020). Recommendation System Using Machine Learning. In *Proceedings of the International Conference on Recent Advances in Computational Techniques (IC-RACT) 2020*. SSRN.
<https://ssrn.com/abstract=3702439> 10
- Meta Learning Recommendation System for Classification. (2020). Air Force Institute of Technology.
<https://scholar.afit.edu/cgi/viewcontent.cgi?article=4630&context=etd> 4
- Metalearning - A Tutorial. (2008). International Conference on Machine Learning and Applications. <https://www.icmla-conference.org/icmla08/slides2.pdf>